

NCBI Molecular Biology Resources

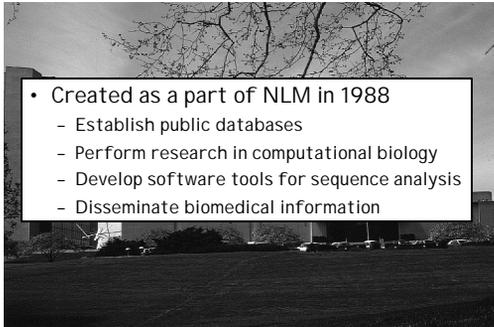
American Society for Microbiology

Peter Cooper
National Center for Biotechnology Information

May 16, 2003

NCBI

The National Center for Biotechnology Information



- Created as a part of NLM in 1988
 - Establish public databases
 - Perform research in computational biology
 - Develop software tools for sequence analysis
 - Disseminate biomedical information

NCBI

Molecular Databases

- Primary
 - Data provided by experimentalist
 - Record maintained by submitter
 - GenBank
- Derivative
 - Value added to primary data
 - compilation
 - curation
 - assembly
 - Record maintained by database staff
 - NCBI Reference Sequences (RefSeq)
 - Protein Data

NCBI

GenBank: NCBI's Primary Sequence Database

Release 135 April 2003

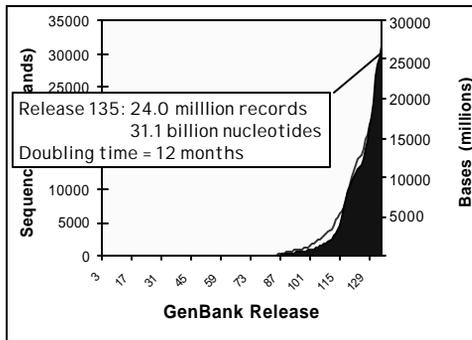
24,027,936	Records
31,099,264,455	Nucleotides
120,000 +	Species

- full release every two months
- incremental and cumulative updates daily
- available only through internet

<ftp://ftp.ncbi.nih.gov/genbank/>

114 Gigabytes

The Growth of GenBank



Traditional GenBank Divisions

- Direct Submissions (Sequin and BankIt)
- Accurate
- Well characterized

BCT	Bacterial and Archeal
INV	Invertebrate
MAM	Mammalian (ex. ROD and PRI)
PHG	Phage
PLN	Plant and Fungal
PRI	Primate
ROD	Rodent
SYN	Synthetic (vectors, synth. genes)
VRL	Viral
VRT	Other Vertebrate

Traditional GenBank Records

LOCUS AF105152 1584 bp mRNA linear VRT 11-FEB-1999
 DEFINITION Danio rerio rhodopsin mRNA, complete cds.
 ACCESSION AF105152
 VERSION AF105152.1 GI:4262520
 KEYWORDS
 SOURCE Danio rerio (zebrafish)
 ORGANISM Danio rerio

LOCUS AF109368 1581 bp mRNA linear VRT 09-AUG-2001
 DEFINITION Danio rerio rhodopsin mRNA, complete cds.
 ACCESSION AF109368
 VERSION AF109368.1 GI:4581736
 KEYWORDS
 SOURCE Danio rerio (zebrafish)
 ORGANISM Danio rerio
 BUKARYOTA; METAZOA; CHORDATA; CRANIATA; VERTEBRATA; EUTELEOSTOMI;
 ACTINOPTERYGII; NEOPTERYGII; TELEOSTEI; OSTARIOPHYSI;
 CYPRINIFORMES; CYPRINIDAE; DANIO.
 1 (bases 1 to 1581)
 AUTHORS Vihalec, T.S., Doro, C.J. and Hyde, D.R.
 TITLE Cloning and characterization of six zebrafish photoreceptor opsin
 cDNAs and immunolocalization of their corresponding proteins
 JOURNAL Vis. Neurosci. 16 (3), 571-585 (1999)
 MEDLINE 9927479
 PUBMED 10549776
 REFERENCE 2 (bases 1 to 1581)
 AUTHORS Vihalec, T.S., Doro, C.J. and Hyde, D.R.
 TITLE Direct Submission
 JOURNAL Submitted (27-NOV-1998) Biological Sciences, University of Notre
 Dame, 264 Galvin Life Science Bldg., Notre Dame, IN 46556, USA

- Archival
- Redundant

IBDN

Bulk GenBank Divisions

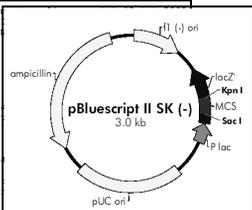
- Batch Submission and htg (email and ftp)
- Inaccurate
- Poorly Characterized

EST Expressed Sequence Tag
STS Sequence Tagged Site
GSS Genome Survey Sequence
HTG High Throughput Genomic

IBDN

Bulk GenBank Records: ESTs

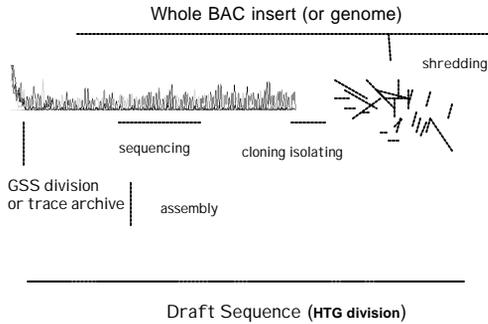
LOCUS BG738754 530 bp
 DEFINITION fp62408.y1 Zebrafish adult retina cDNA
 IMAGE:4796599 5' similar to FR:Q9P999
 ACCESSION BG738754
 VERSION BG738754.1 GI:14088443
 KEYWORDS EST
 SOURCE Danio rerio (zebrafish)
 ORGANISM Danio rerio
 BUKARYOTA; METAZOA; CHORDATA; CRANIATA; VERTEBRATA; EUTELEOSTOMI;
 ACTINOPTERYGII; NEOPTERYGII; TELEOSTEI; OSTARIOPHYSI;
 CYPRINIFORMES; CYPRINIDAE; DANIO.
 1 (bases 1 to 530)
 AUTHORS Clark, N., Johnson, S.L., Lehrach, H., Miller, L., Kucaba, T., Martin, J., Tegetmeier, M., Theising, B., Allen Swaller, T., Gibbons, M., Pape, D., Ha Kohn, S., Shin, T., Jackson, V., Carde and Wilson, R.
 TITLE WashU Zebrafish EST Project 1998
 JOURNAL Unpublished (1998)
 COMMENT Contact: Stephen L. Johnson
 Washington University School of
 4444 Forest Park Parkway, Box 85
 Tel: 314 286 1800
 Fax: 314 286 1810
 Email: zbrfish@watson.wustl.edu
 Library constructed by: Chandra
 Sequencing by: Washington University Genome Sequencing Center Clone
 distribution: Ressourcenzentrum PrimatDatenbank, Berlin, Germany
 (web address: www.rpd.de)
 Seq primer: T3 ET from Amersham
 High quality sequence stop: 423.



- first pass, single read cDNA
- largest GenBank division
- gbdiv_est[Properties]

IBDN

Genome Sequencing - HTG, GSS, (WGS)



Bulk GenBank Records: GSS

```

LOCUS      M247846                551 bp    DNA             linear   08F 14-09-03
DEFINITION Danio rerio genomic clone BREV-23509, genomic survey sequence.
ACCESSION  M247846
VERSION   M247846.1  GI:28170522
KEYWORDS  GSS.
SOURCE    Danio rerio (zebrafish)
  ORGANISM  Danio rerio
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Actinopterygii; Neopterygii; Teleostei; Ostariophysi;
            Cypriniformes; Cyprinidae; Danio.
REFERENCE  1 (bases 1 to 551)
  AUTHORS  Humphrey,S.J., Huckle,E. and Durhan,J.L.
  TITLE    Direct Subcloning
  JOURNAL  Submitted (13-MAR-2003) The Sanger Institute, Wellcome Trust Genome
  Campus, Hinxton, Cambridgeshire, CB1 1BA, UK. E-mail enquiries
  humphrey@sanger.ac.uk (humphreys)
  COMMENT  This sequence was generated from the 896 end of BAC 23509, 239
  part of the Danio/BAC library created by S. Harcourt and H.
  Koyama. Further details:
  http://www.sanger.ac.uk/projects/D_rerio/
  FEATURES     source
            location/Qualifiers
            1..551
            /organism="Danio rerio"
            /mol_type="genomic DNA"
            /db_xref="taxon:7955"
            /clone="BREV-23509"
            /issue_type="Genis"
            /note="Vector: pBluescript-2.1"
  BASE COUNT  205 a   71 c   110 g   165 t
  ORIGIN
1  cttatttaa attttctat gttgaaaca aaacagcaa tt
61  aactcaagc acttaactg tgaagagg ttataaaga at
121  acttgagac tgaactatg gagaactcg aaataccaa at
181  gctgggtgc ctatgaca aaattgac aaacagga at
241  ggtcattaa ttaacatgc agtggagca acaccaag at
301  aatggtgca atagagaa aggggttgc atgggaa at
361  aactggaaa ttctgtggt atcaaatct taatttcca at
421  aactgcgca atgactcaga cttctgata aagataaa gbaattttc aatcttgc
481  tgaagagca aacttcaga atctcttaa agagactaa gaactcacc aggaactgc
541  aggtataaa t
  //
  
```

- first pass, single read gDNA
- surveys of BAC libraries
- `gbdiv_gss[Properties]`
- zebrafish GSS: 159, 024

HTG Division: High Throughput Genome

Zebrafish PAC Clone

phase 1 → ← → → ← HTG
Acc = AC109580.1

phase 2 → → → → HTG
Acc = AC109580.11

phase 3 → → → → VRT
Acc = AC109580.14

`gbdiv_htg[Properties]`
zebrafish: 2,958

40,000 to > 350,000 bp

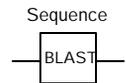
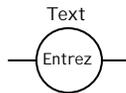
RefSeq: NCBI's Derivative Sequence Database

- **Curated transcripts and proteins [NM_, NP_]**
 - reviewed
 - human, mouse, rat, cow, fruit fly, zebrafish, arabidopsis, *C.elegans*
- **Model transcripts and proteins [XM_, XP_]**
- **Assembled Genomic Regions [NT_, NW_]**
 - draft human genome
 - mouse genome
- **Chromosome records [NC_]**
 - microbial
 - organelle

`srcdb_refseq[Properties]`

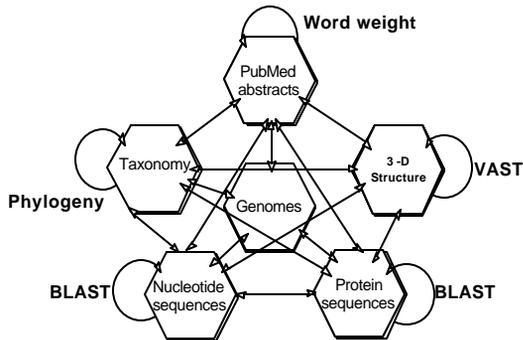
IBDN

Web Access



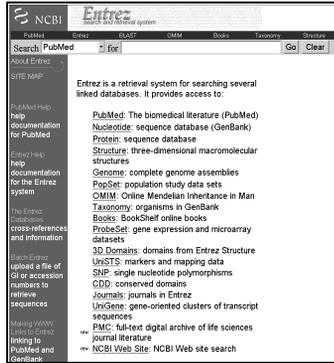
IBDN

Entrez: Database Integration



IBDN

WWW Entrez



Sequence Similarity Searching

Basic Local Alignment Search Tool

Why do we need similarity searching?

- ◆ Identification and annotation
 - Incomplete or no annotations (GenBank)
 - Incorrectly annotated sequences
- ◆ Evolutionary relationships
 - homologous molecules **may** have similar functions

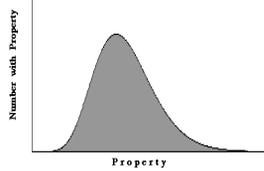
Basic Local Alignment Search Tool

- Widely used similarity search tool
- Heuristic approach based on Smith Waterman algorithm
- Finds best local alignments
- Provides statistical significance
- All combinations (DNA/Protein) query and database.
 - DNA vs DNA
 - DNA translation vs Protein
 - Protein vs Protein
 - Protein vs DNA translation
 - DNA translation vs DNA translation
- www, standalone, and network clients

IBSON

Local Alignment Statistics

High scores of local alignments between two random sequences follow Extreme Value Distribution



For ungapped alignments:
 Expected number with score S or greater
 $E = Kmne^{-lS}$
 or
 $E = mn2^{-S}$

K = scale for search space
 l = scale for scoring system
 $S = \text{bitscore} = (lS - \ln K) / \ln 2$

<http://www.ncbi.nlm.nih.gov/BLAST/tutorial/Aitschul-1.html>

IBSON

Scoring Systems

•Position Independent Matrices

•Nucleic Acids – identity matrix

•Proteins

- PAM Matrices (Percent Accepted Mutation)
 - Implicit model of evolution
 - Higher PAM number all calculated from PAM1
 - PAM250 widely used
- BLOSUM Matrices (BLOck SUBstitution Matrices)
 - Empirically determined from alignment of conserved blocks
 - Each includes information up to a certain level of identity
 - BLOSUM62 widely used

•Position Specific Score Matrices (PSSMs)

- PSI and RPS BLAST

	A	G	C	T
A	+1	-3	-3	-3
G	-3	+1	-3	-3
C	-3	-3	+1	-3
T	-3	-3	-3	+1

IBSON

BLAST Databases: Non-redundant protein

nr (non-redundant protein sequences)

- GenBank CDS translations
- NP_ RefSeqs
- Outside Protein
 - PIR, **Swiss-Prot**, PRF
- **PDB** (sequences from structures)

BLAST Databases: Nucleic Acid

- **nr (nt)**
 - Traditional GenBank Divisions
 - NM_ and XM_ RefSeqs
- **dbest**
 - EST Division
- **htgs**
 - HTG division
- **gss**
 - GSS division
- **chromosome**
 - NC_ RefSeqs

Service Addresses

- **General Help** info@ncbi.nlm.nih.gov
- **BLAST** blast-help@ncbi.nlm.nih.gov

NCBI

NCBI
**Field
Guide**

NCBI Training

A Field Guide to NCBI Resources

To host a course at your institution, write to Peter Cooper.

Course Description

The National Center for Biotechnology Information (NCBI) provides a **Field Guide to GenBank and NCBI Molecular Biology Resources** as lecture and hands-on computer workbooks on GenBank and related databases covering effective use of the Entrez database and search services for BLAST, similarity search engines and genome data and related resources.

Now featuring the NCBI assembly of the draft human genome, the NCBI's **mouse whole genome shotgun assembly**, the updated human and mouse map viewers <http://www.ncbi.nlm.nih.gov/blast/>.

Day Long Training Course
Audience: Research Biologists
 • 3-hour Lecture
 • 2-hour hands on practice
 • One-on-one consults

Course Objectives

- Identify key resources
- Use BLAST/BLAT/BLASTZ
- Complete Molecular Sequence in Entrez
- Multiple Sequence Alignment
- Simple protein models
- Clustal
- GenBank files (FASTA)
- The NCBI Tree Browser/Phylo
- The Map Viewer
- Mouse and Other Genomes

2003 Courses

<p>Mar 6 Utah Valley State College Orms</p> <p>Mar 10 Virginia Commonwealth University Bullinswood</p> <p>Mar 17 Heidelberg, Germany EZZD</p> <p>Mar 23 Texas A & M University College Station, Texas</p> <p>Mar 27 National Library of Medicine Bethesda, Maryland</p> <p>Apr 9 Frederick Hutchinson Cancer Research Center Seattle, Washington</p> <p>Mineral Cancer Institute Frederick, Maryland prosser@college</p> <p>Williamsport, Pennsylvania</p> <p>North Carolina A & T State University Reno/Raleigh</p> <p>Iowa State University Agricultural Campus Weaver</p> <p>Iowa State University Iowa</p> <p>May 9 Life Journal Symposium on Pollutant Exposures in Marine Organisms (FRMO) Safety Harbor, Florida</p> <p>May 16 National "Wiring" Open Bioinformatics Tools and Databases in the Undergraduate Biology Curriculum American Society of Microbiology Undergraduate Microbiology Education Conference University of Maryland College Park</p>	<p>Janina Price sd</p> <p>Catherine pc</p> <p>Catherine rm</p> <p>pc</p> <p>Robert Sewell wsm</p> <p>es</p> <p>Peter Cooper pc</p> <p>Barry Stoddard sd</p> <p>rm</p> <p>Robert Mackley pc</p> <p>Jeff Newman wsm</p> <p>es</p> <p>Melanie Ott sd</p> <p>Wolke rm</p> <p>Jan Holford tha</p> <p>Peter Sheridan tha</p> <p>pc</p> <p>pc</p>
----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

NCBI

10